# Speedrunning the Lakehouse
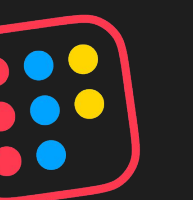
A composable FaaS over object storage

# Ciao, I'm Jacopo!

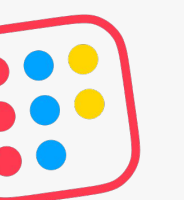| Co-founder and CTO at **Bauplan**.
Backed by IE, SPC, Wes McKinney, Spencer Kimball, Chris Re et al..

| Started the "Reasonable Scale" movement.
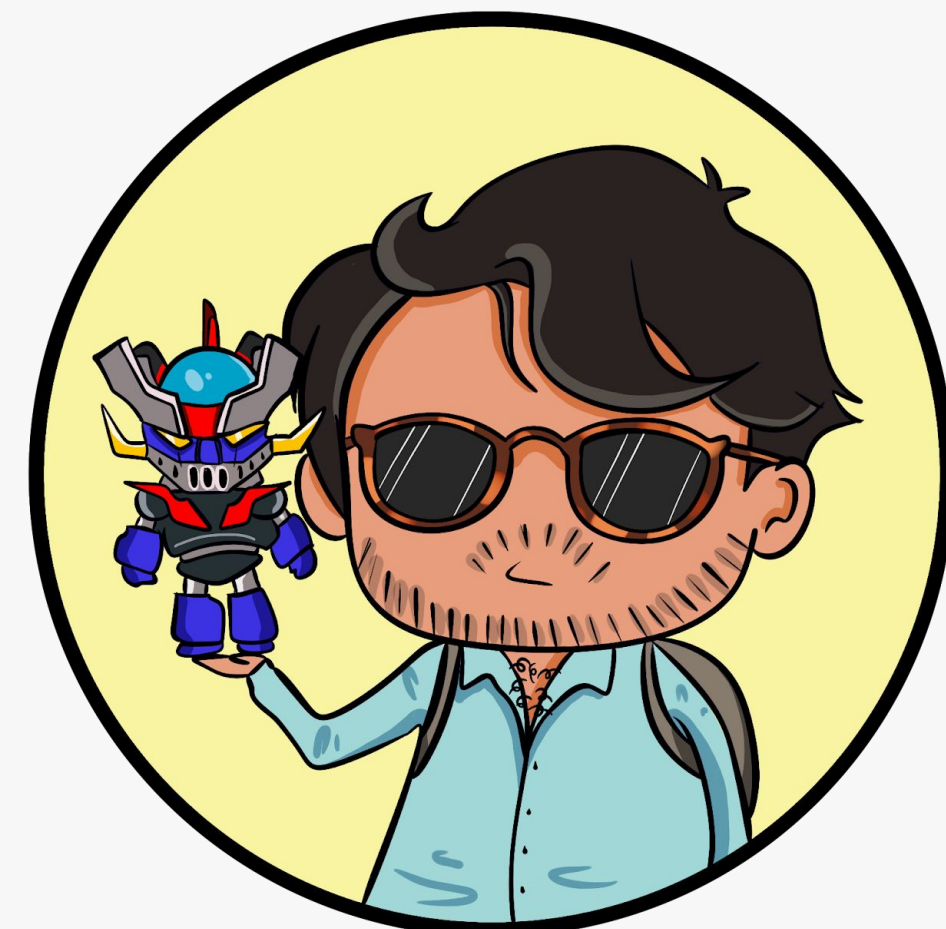Co-founder at Tooso and lead AI at TSX:CVO after the acquisition.

| 10 years up and down the stack in R&D, product, open source
ICML, KDD, VLDB, NAACL, SIGIR, WWW et al.., >2k stars, >50M+ downloads.

# It takes a (distributed) village

| Matt, Ciro, Luca, Nate, Vlad (and others, unfortunately without a chibi) share with me the credit for whatever value these ideas may have.
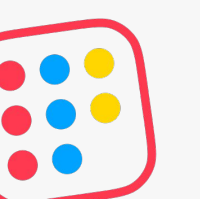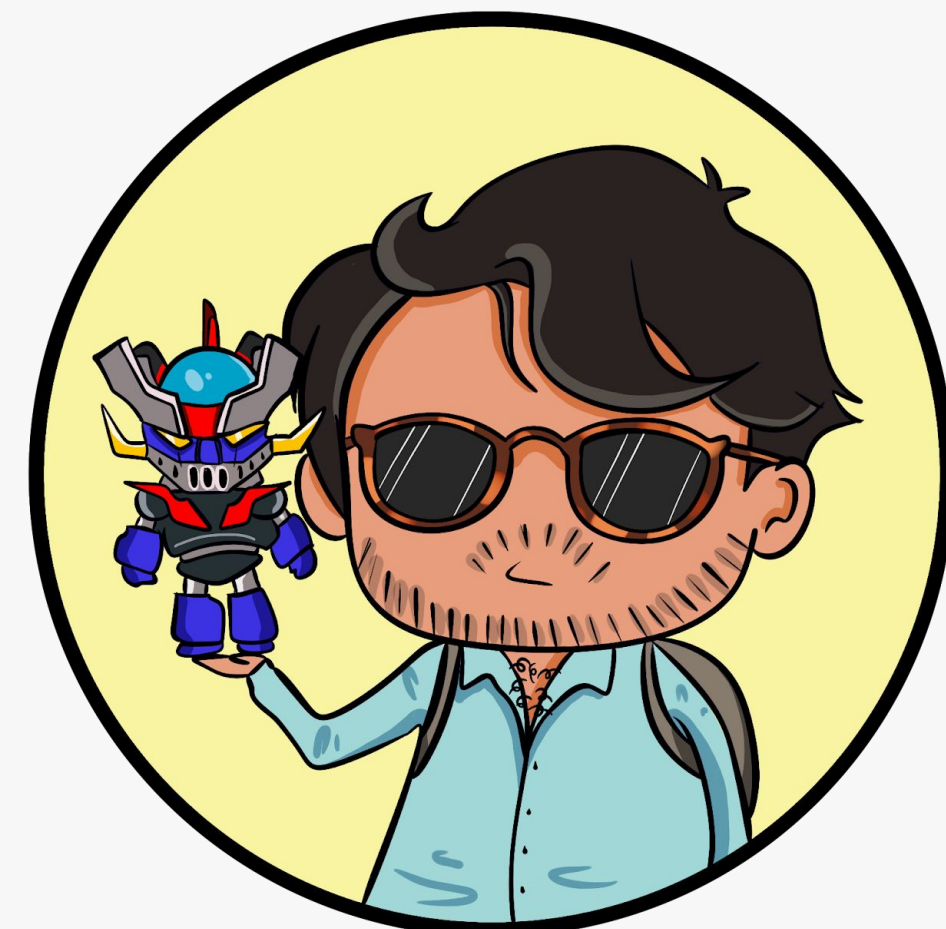
# It takes a (distributed) village

| Matt, Ciro, Luca, Nate, Vlad (and others, unfortunately without a chibi) share with me the credit for whatever value these ideas may have.
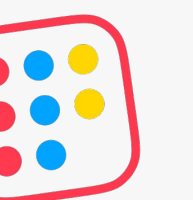
| Obviously, all the remaining mistakes are theirs 😁

"bauplan is [a system] fully built using composable principles (...). It is refreshing to learn about a real–life system built using such architectural principles."
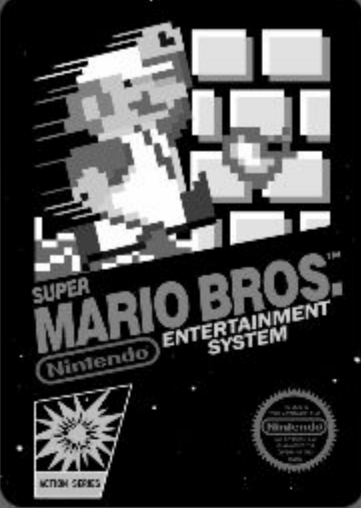
*Reviewer #2*

speed·run
/ˈspēdˌrən/

verb

gerund or present participle: **speedrunning**

complete (a video game, or level of a game) as fast as possible.

"I used to be able to speedrun this game in less than 20 minutes"

**Super Mario Bros.** [1985]
Super Mario Series
NES  SNES  WiiVC  +14

Category extensions     Discord

Leaderboards   News 8   Guides 42   Resources 44

Any%   Warpless   Any% All-Stars   Warpless All-Stars

Version
NTSC  PAL

Filters          Show rules

| # | Player | | Time |
|---|--------|---|------|
| ☆ | Niftski | | 4m 54s 565ms |
| | averge11 | | 4m 54s 748ms |
| | Tree_05 | | 4m 54s 864ms |

**Time**

4m 54s 565ms

4m 54s 748ms

4m 54s 864ms

# Speedrunning a Lakehouse? Really?

Dremio
2B USD

Snowflake
66B

Databricks
100B

Fabric
???B

# Speedrunning a Lakehouse? Really?

Bauplan

Dremio

Snowflake

Databricks

Fabric

# 1. Simplicity
# 2. Composability

# 🔲 Mo (mental) models, mo problems

| Interaction | UX | Infrastructure |
| --- | --- | --- |
| **Traditional DLH** | | |
| Batch pipeline | Submit API | One-off cluster |
| Dev. pipeline | Notebook Session | Dev. cluster |
| Inter. query | Web Editor (JDBC Driver) | Warehouse |

### *Eudoxia*: a FaaS scheduling simulator for the composable lakehouse

Tapan Srivastava*
tapansriv@uchicago.edu
University of Chicago
Chicago, Illinois, USA

Jacopo Tagliabue*
jacopo.tagliabue@bauplanlabs.com
Bauplan Labs
New York, USA

Ciro Greco
ciro.greco@bauplanlabs.com
Bauplan Labs
New York, USA

**ABSTRACT**

Due to the variety of its target use cases and the large API surface area to cover, a data lakehouse (DLH) is a natural candidate for a composable data system. *Bauplan* is a composable DLH built on "spare data parts" and a unified Function-as-a-Service (FaaS) runtime for SQL queries and Python pipelines. While FaaS simplifies both building and using the system, it introduces novel challenges in scheduling and optimization of data workloads. In this work, starting from the programming model of the composable DLH

data lake and warehouse, such as cheap and durable foundation through object storage, compute decoupling, multi-language support, unified table semantics, and governance [19].

The breadth of DLH use cases makes it a natural target for the philosophy of composable data systems [23]. In this spirit, *Bauplan* is a DLH built from "spare parts" [31]: while presenting to users a unified API for assets and compute [30], the system is built from modularized components that reuse existing data tools through novel interfaces: e.g. Arrow fragments for differential caching [29],

May 2025

```
pip install bauplan
bauplan checkout my-branch
bauplan run
```

"Simplex Sigillum Veri"

# A sample pipeline

**transactions**

| ID | USD | COUNTRY |
|----|-----|---------|
| 13 | 44 | US |
| 144 | 13 | IT |
| 146 | 1 | IT |

```python
def euro_selection(
    df=transactions
):
    _df =
transform_input(df)
    return _df
```

**euro_selection**

| ID | USD | COUNTRY |
|----|-----|---------|
| 144 | 13 | IT |
| 146 | 1 | IT |

```python
def usd_by_country(
    df=euro_selection
):
    _df =
transform_input(df)
    return _df
```

**usd_by_country**

| COUNTRY | USD |
|---------|-----|
| IT | 14 |

# High-level view of a bauplan run

# ⬚ High-level view of a `bauplan` run



USER

BAUPLAN AWS

CUSTOMER AWS

WORKER #2

(4)

PUBLIC APIs —(3)

WORKER #1

(2)

SDK —(1)

DATA CATALOG

FILES

CLI / SDK

Control Plane

Data Plane

# Pipelines are chained functions (Batch / Dev)



functions

**user *f1***

```python
def euro_selection(
    df=euro_selection
):

    _df = transform_input(df)
    return _df
```

**user *f2***

```python
def usd_by_country(
    df=euro_selection
):
    _df = transform_input(df)
    return _df
```

S3 read *f*

S3 write *f*

APACHE ARROW

functions

# Queries are chained functions as well!

# Everything is a function, or "OnlyFaas"

| Easy to reason about

  - Simple abstractions, "looks like code"
  - A unified compute model, "everything is a function"



VLDB 2023: Building a serverless Data Lakehouse from spare parts

# PROs: one mental model to rule them all

| Can we re-use existing FaaS? **NO!!!**

- Resource limitations

- No "DAG awareness"

- Slow feedback loop

| Interaction | UX | Infrastructure |
|---|---|---|
| Traditional DLH | | |
| Batch pipeline | Submit API | One-off cluster |
| Dev. pipeline | Notebook Session | Dev. cluster |
| Inter. query | Web Editor (JDBC Driver) | Warehouse |

# CONs: we need a new a FaaS-for-data

| New programming model

   – Express data and code dependencies

| New runtime

   – Function lifecycle

   – Scheduling

| Interaction | UX | Infrastructure |
|---|---|---|
| **Traditional DLH** | | |
| Batch pipeline | Submit API | One-off cluster |
| Dev. pipeline | Notebook Session | Dev. cluster |
| Inter. query | Web Editor (JDBC Driver) | Warehouse |

# New programming model

clear "division of labor"
between platform and users

# New programming model

bau.py

```python
@bauplan.model()
@bauplan.python(
    "3.11",
    pip={"polars": "1.33.0"}
)
def euro_selection(
    data=bauplan.Model(
        "transactions",
        columns=["id", "usd", "country"],
        filter="eventTime BETWEEN 2023-01-01 AND 2023-02-01"
    )
):
    # filtering here
    # return a dataframe
    return _df
```

bau.py

```python
@bauplan.model(materialize=True)
@bauplan.python(
    "3.10",
    pip={"polars": "0.8.8"}
)
def usd_by_country(
    data=bauplan.Model("euro_selection")
):
    # aggregation here
    # return a dataframe
    return _df
```

# New programming model

**User code here!**

bau.py

```python
@bauplan.model()
@bauplan.python(
    "3.11",
    pip={"polars": "1.33.0"}
)
def euro_selection(
    data=bauplan.Model(
        "transactions",
        columns=["id", "usd", "country"],
        filter="eventTime BETWEEN 2023-01-01 AND
2023-02-01"
    )
):
    # filtering here
    # return a dataframe
    return _df
```

bau.py

```python
@bauplan.model(materialize=True)
@bauplan.python(
    "3.10",
    pip={"polars": "0.8.8"}
)
def usd_by_country(
    data=bauplan.Model("euro_selection")
):
    # aggregation here
    # return a dataframe
    return _df
```

# New programming model

**Signature Table(s)->Table**

bau.py

```python
@bauplan.model()
@bauplan.python(
    "3.11",
    pip={"polars": "1.33.0"}
)
def euro_selection(
    data=bauplan.Model(
        "transactions",
        columns=["id", "usd", "country"],
        filter="eventTime BETWEEN 2023-01-01 AND 2023-02-01"
    )
):
    # filtering here
    # return a dataframe
    return _df
```

bau.py

```python
@bauplan.model(materialize=True)
@bauplan.python(
    "3.10",
    pip={"polars": "0.8.8"}
)
def usd_by_country(
    data=bauplan.Model("euro_selection")
):
    # aggregation here
    # return a dataframe
    return _df
```

# New programming model

**Infra-as-code**

bau.py

```python
@bauplan.model()
@bauplan.python(
    "3.11",
    pip={"polars": "1.33.0"}
)
def euro_selection(
    data=bauplan.Model(
        "transactions",
        columns=["id", "usd", "country"],
        filter="eventTime BETWEEN 2023-01-01 AND
2023-02-01"
    )
):
    # filtering here
    # return a dataframe
    return _df
```

bau.py

```python
@bauplan.model(materialize=True)
@bauplan.python(
    "3.10",
    pip={"polars": "0.8.8"}
)
def usd_by_country(
    data=bauplan.Model("euro_selection")
):
    # aggregation here
    # return a dataframe
    return _df
```

# New programming model

**I/O chaining**

bau.py

```python
@bauplan.model()
@bauplan.python(
    "3.11",
    pip={"polars": "1.33.0"}
)
def euro_selection(
    data=bauplan.Model(
        "transactions",
        columns=["id", "usd", "country"],
        filter="eventTime BETWEEN 2023-01-01 AND
2023-02-01"
    )
):
    # filtering here
    # return a dataframe
    return _df
```
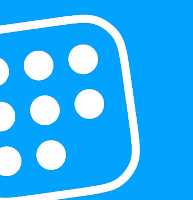
bau.py

```python
@bauplan.model(materialize=True)
@bauplan.python(
    "3.10",
    pip={"polars": "0.8.8"}
)
def usd_by_country(
    data=bauplan.Model("euro_selection")
):
    # aggregation here
    # return a dataframe
    return _df
```
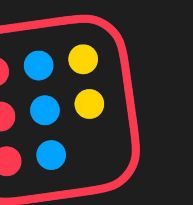
# New runtime

we can't just "run user functions", which
is a challenge and opportunity

**bauplan run** =
$$\begin{bmatrix} \text{plan} \\ + \\ \text{environment} \\ + \\ \text{data movement} \end{bmatrix}$$

# ⊞ Planning

**USER CODE**

**PLATFORM CODE**

bau.py

```python
@bauplan.model()
@bauplan.python(
    "3.11",
    pip={"polars": "0.8.8"}
)
def euro_selection(
    data=bauplan.Model(
        "transactions",
        columns=["id", "usd", "country"],
        filter="eventTime BETWEEN 2023-01-01 AND
2023-02-01"
    )
):
    # filtering here
    # return a dataframe
    return _df
```
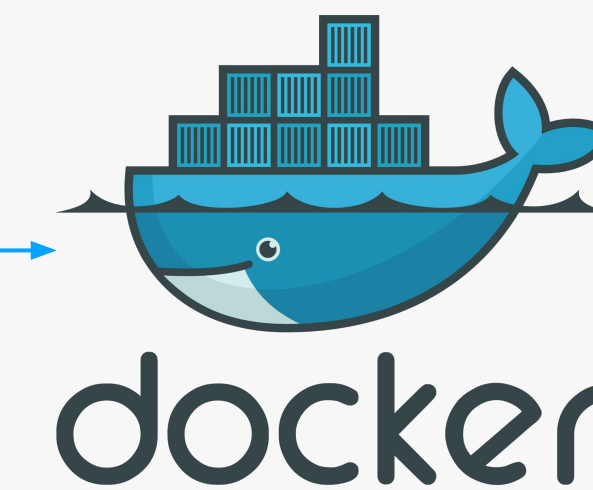
```
RUN ...
…
```

```
obj.get(Range='bytes=32-64')['Body']
…
```

# Planning

Logical

transactions → **f** → euro_selection → **g** → usd_by_country

Physical

Worker

cache

cont.
factory

APACHE **ARROW**

APACHE **ARROW**

CDMS@VLDB2025: Dag lakehouse planning with an ephemeral and embedded graph database

# Environment

```
@bauplan.python(
    "3.11",
    pip={"polars": "0.8.8"}
)
```

**planner**

Dependency graph

Polars 0.8

Package 1.0

Package 2.1

Package 2.0

bauplan cloud

# Environment

```python
@bauplan.python(
    "3.11",
    pip={"polars": "0.8.8"}
)
```

planner

bauplan cloud

Dependency graph

Polars 0.8

Package 1.0

Package 2.1

Package 2.0

worker

customer cloud

Install to ...

Polars 0.8

Package 1.0

Package 2.1

Package 2.0

# Environment

```python
@bauplan.python(
    "3.11",
    pip={"polars": "0.8.8"}
)
```

**planner**

Dependency graph

Polars 0.8

Package 1.0

Package 2.1

Package 2.0

**bauplan cloud**

**worker**

Install to ...

Polars 0.8

Package 1.0

Package 2.1

Package 2.0

mounted packages

user code

**customer cloud**

# ⠿ Environment: assemble, don't build

| **NO** Docker, **NO** bandwidth bottlenecks, **NO** ECR update

| Functions are ephemeral: no lifecycle management.

| **Adding a package is 15× faster than AWS Lambda**

## Table 2: Time to add *Prophet* to a serverless DAG

| Task | Seconds |
|------|---------|
| **AWS Lambda**[4] | |
| Update ECR container and function | 130 (80 + 50) |
| **Snowpark** | |
| Update Snowpark container | 35 |
| *bauplan* | |
| Update runtime | 5 / 0 (cache) |

# Data movement: Arrow everywhere + zero-copy

```
data=bauplan.Model(
    "transactions",
    columns=["id", "usd", "country"],
    filter="..."
)
```

APACHE
**ARROW** 》》》

| Across workers, an Arrow stream is as fast as local parquet files (**B**)

| Within a worker, tables can be zero-copy shared between functions (**C**)



Figure 1: **Communication: Degrees of Zero Copy**

# 🎲 Data movement: Arrow everywhere + zero-copy

**Table 3: Reading a dataframe from a parent (*c5.9xlarge*), avg. (SD) over 5 trials**

|  | *10M rows (6 GB)* | *50M rows (30 GB)* |
|---|---|---|
| Parquet file in S3 | 1.26 (0.14) | 6.14 (0.98) |
| Parquet file on SSD | 0.92 (0.09) | 4.37 (0.15) |
| Arrow Flight | 0.96 (0.01) | 4.69 (0.01) |
| Arrow IPC | **0.01 (0.00)** | **0.03 (0.01)** |

**RQ:** *is D even feasible?*

**Zerrow: True Zero-Copy Arrow Pipelines in *Bauplan***

Yifan Dai*, Jacopo Tagliabue*, Andrea Arpaci-Dusseau*,
Remzi Arpaci-Dusseau*, Tyler R. Caraza-Harter**
* *University of Wisconsin–Madison*, * *Bauplan Labs*

*Abstract.* Bauplan is a FaaS-based lakehouse specifically built for data pipelines: its execution engine uses Apache Arrow for data passing between the nodes in the DAG. While Arrow is known as the "zero copy format", in practice, limited Linux kernel support for shared memory makes it difficult to avoid copying entirely. In *this* work, we introduce several new techniques to eliminate nearly all copying from pipelines: in particular, we implement a new kernel module that performs de-anonymization, thus eliminating a copy to intermediate data. We conclude by sharing our preliminary evaluation on different workloads types, as well as discussing our plan for future improvements.
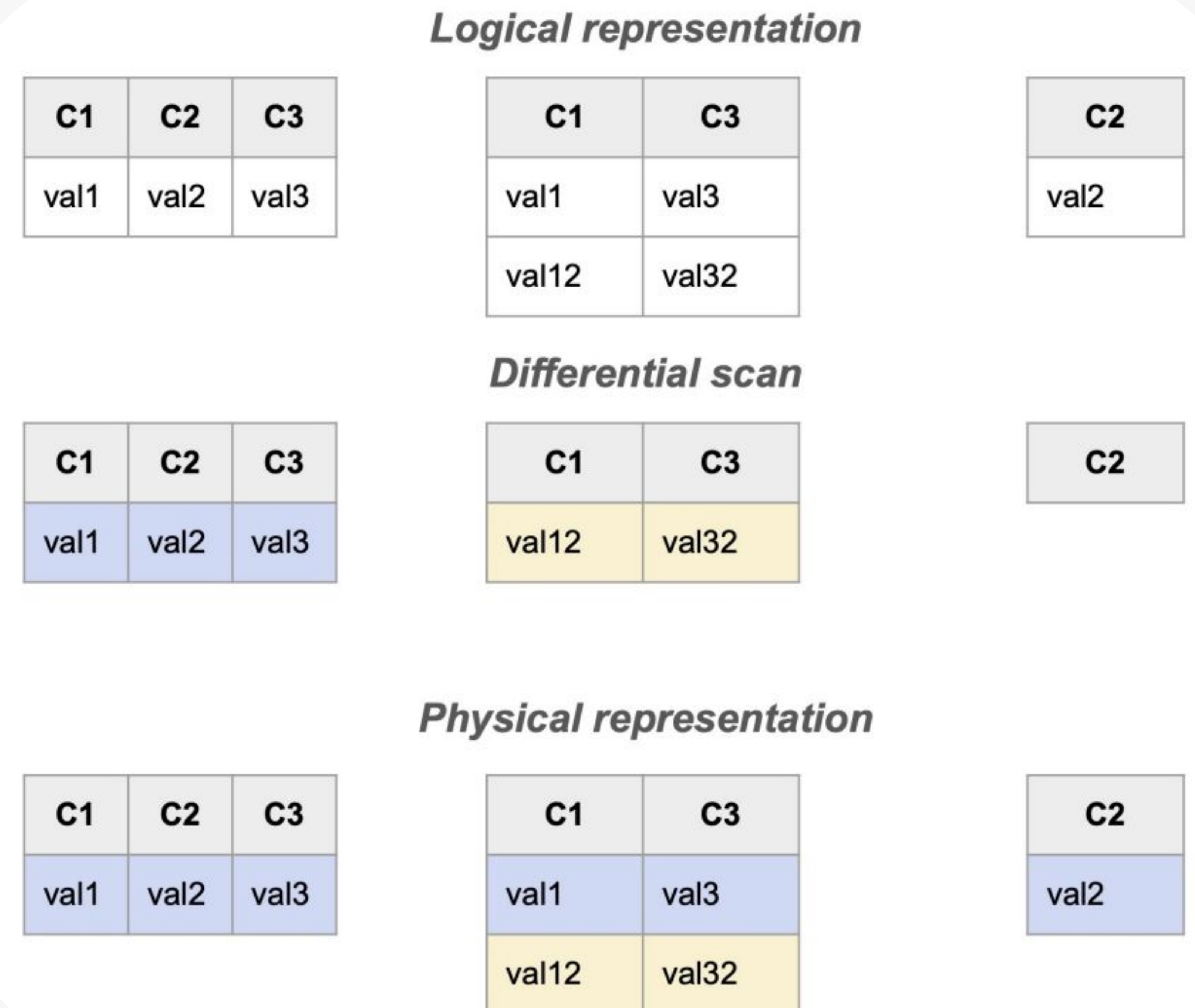
be mapped by multiple downstream nodes. Unfortunately, simply using Arrow for inter-node communication does not eliminate several sources of copying and duplication in data pipelines. First, many tools and libraries that return Arrow data allocate space with malloc, which uses anonymously mapped memory without a backing file; operating systems (including Linux) do not typically support sharing of anonymous memory, so unless all libraries in the Arrow ecosystem are rewritten to use shared memory, a copy to shared memory is necessary. Second, DAG nodes must perform copies when Arrow output overlaps with Arrow input (*e.g.*, the node adds a column to an input table), as the existing Arrow IPC protocol does not provide a way to identify or reference such overlap. Finally, when independent DAGs deserialize the same data from on-disk formats (*e.g.*, Parquet files) to Ar-

**1 Introduction**

Data pipelines are a popular programming paradigm for data

# Scans do not repeat themselves, but they often rhyme

**Differential cache:**

|     U1: "SELECT c1, c2, c3 FROM t WHERE eventTime BETWEEN 2023–01–01 AND 2023–02–01"

|     U2: "SELECT c1, c3 … BETWEEN 2023–01–01 AND 2023–03–01"

|     U1: "SELECT c2 … BETWEEN 2023–01–01 AND 2023–01–02"



*Logical representation*

| C1 | C2 | C3 |
|------|------|------|
| val1 | val2 | val3 |

| C1 | C3 |
|-------|-------|
| val1 | val3 |
| val12 | val32 |

| C2 |
|------|
| val2 |

*Differential scan*

| C1 | C2 | C3 |
|------|------|------|
| val1 | val2 | val3 |

| C1 | C3 |
|-------|-------|
| val12 | val32 |

| C2 |
|----|
| |

*Physical representation*

| C1 | C2 | C3 |
|------|------|------|
| val1 | val2 | val3 |

| C1 | C3 |
|-------|-------|
| val1 | val3 |
| val12 | val32 |

| C2 |
|------|
| val2 |

# RQ: *how do you manage concurrent functions?*

**Eudoxia: a FaaS scheduling simulator for the composable lakehouse**

Tapan Srivastava*
tapansriv@uchicago.edu
University of Chicago
Chicago, Illinois, USA

Jacopo Tagliabue*
jacopo.tagliabue@bauplanlabs.com
Bauplan Labs
New York, USA

Ciro Greco
ciro.greco@bauplanlabs.com
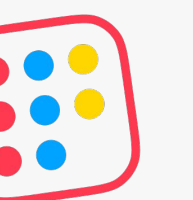Bauplan Labs
New York, USA

## ABSTRACT

Due to the variety of its target use cases and the large API surface area to cover, a data lakehouse (DLH) is a natural candidate for a composable data system. *Bauplan* is a composable DLH built on "spare data parts" and a unified Function-as-a-Service (FaaS) runtime for SQL queries and Python pipelines. While FaaS simplifies both building and using the system, it introduces novel challenges in scheduling and optimization of data workloads. In this work, starting from the programming model of the composable DLH, we characterize the underlying scheduling problem and motivate simulations as an effective tools to iterate on the DLH. We then

data lake and warehouse, such as cheap and durable foundation through object storage, compute decoupling, multi-language support, unified table semantics, and governance [19].

The breadth of DLH use cases makes it a natural target for the philosophy of composable data systems [23]. In this spirit, *Bauplan* is a DLH built from "spare parts" [31]: while presenting to users a unified API for assets and compute [30], the system is built from modularized components that reuse existing data tools through novel interfaces: e.g. Arrow fragments for differential caching [29], Kuzu for DAG planning [18], DuckDB as SQL engine [24], Arrow Flight for client-server communication [6].
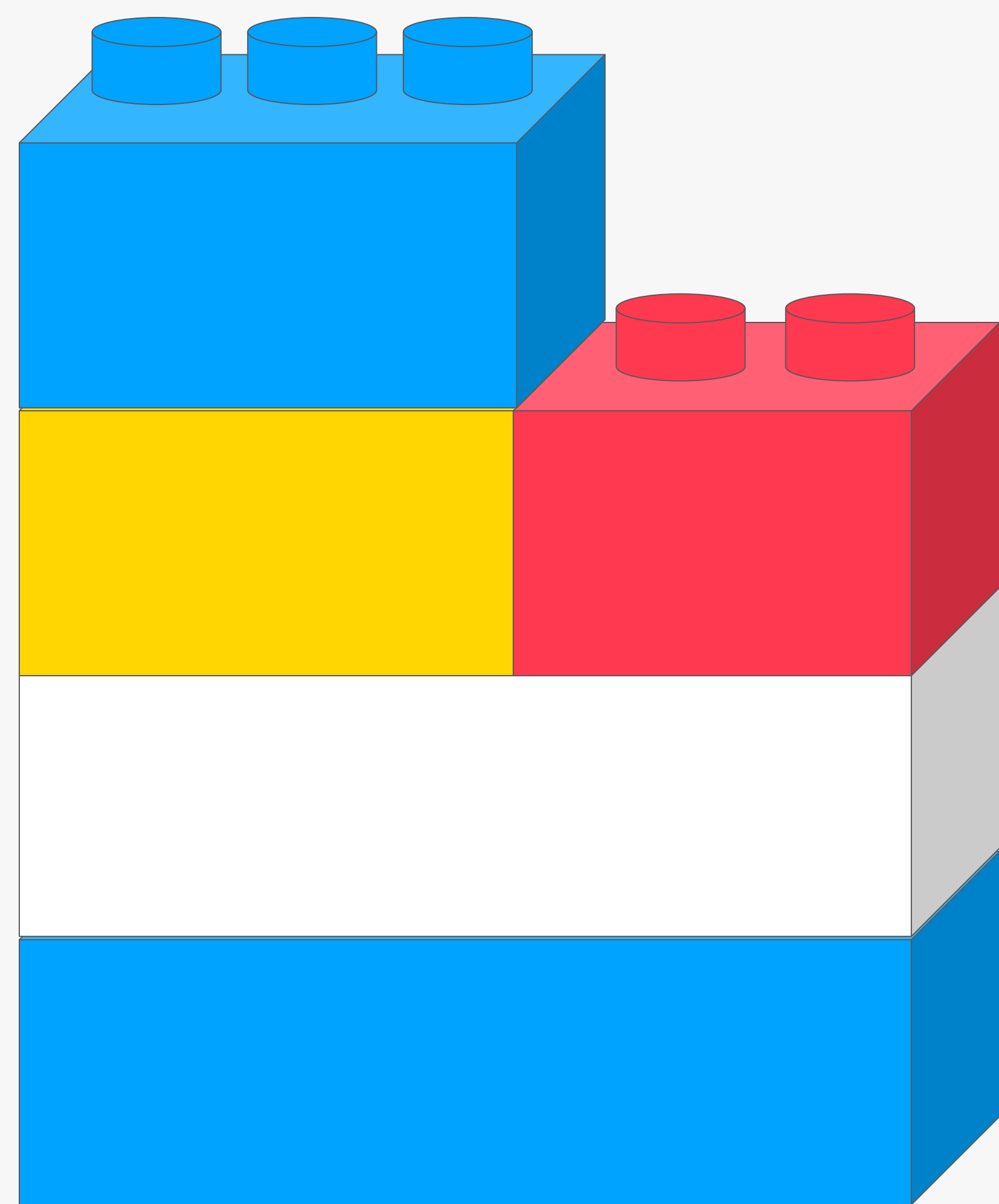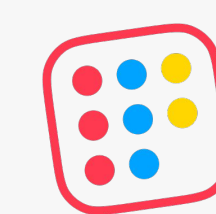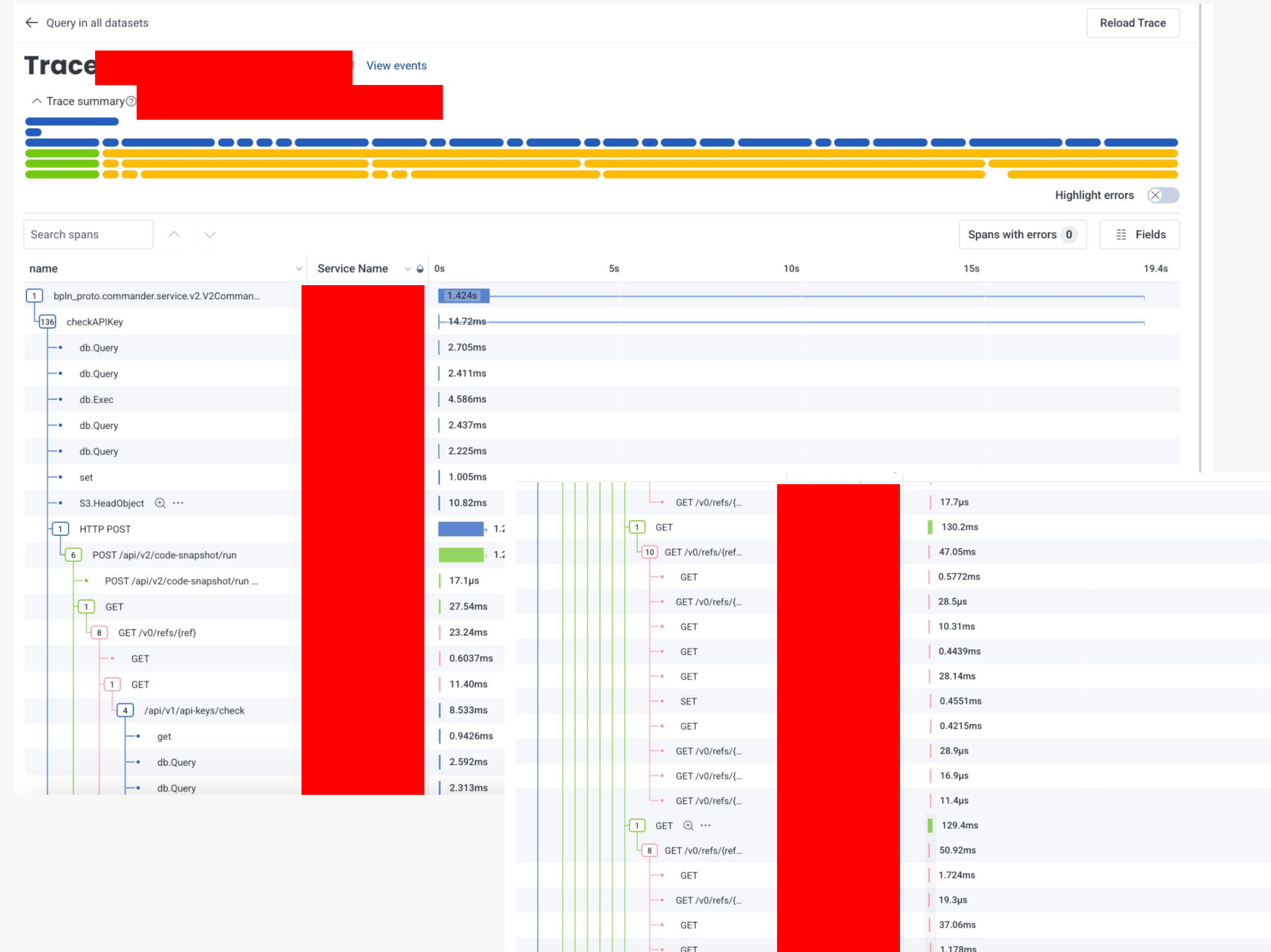
19 May 2025

# The Composable Lakehouse

composable

# Low-level view of a `bauplan` run

A single *run* spans hundreds of traces, across dozens of services, ranging from hyper-scaler PaaS to obscure open-source libraries.

# The *composable* lakehouse

**Bauplan is built around some core competencies plus "spare parts":**

| FaaS runtime and abstractions are new

| Our query engine is a fork of DuckDb

| Our catalog is a fork of Nessie

| Our DAG planner is built on Kuzu

| Our Iceberg client is a fork of PyIceberg

## Building a serverless Data Lakehouse from spare parts*

Jacopo Tagliabue[1,2,*], Ciro Greco[1] and Luca Bigon[1,†]

[1]*Bauplan, New York City, United States*
[2]*Tandon School of Engineering, NYU, New York City, United States*

**Abstract**
The recently proposed Data Lakehouse architecture is built on open file formats, performance, and first-class support for data transformation, BI and data science: while the vision stresses the importance of lowering the barrier for data work, existing implementations often struggle to live up to user expectations. At *Bauplan*, we decided to build a new serverless platform to fulfill the Lakehouse vision. Since building from scratch is a challenge unfit for a startup, we started by re-using (sometimes unconventionally) existing projects, and then investing in improving the areas that would give us the highest marginal gains for the developer experience. In this work, we review user experience, high-level architecture and tooling decisions, and conclude by sharing plans for future development.

**Keywords**
data lakehouse, data pipelines, serverless, reasonable scale, containerized execution

## 1. Introduction

[2] argues that the popular data warehouse architecture will soon be replaced by a new architectural pattern, the Data Lakehouse (DLH). A DLH is built on open file formats (e.g. Parquet), exceptional performance, and first-class support for engineering (data transformation), analytics (BI) and inferential (data science) use cases. The vision of such architecture is first and foremost about flexibility, making it possible for organizations to choose different ways to operationalize data depending on data volumes, use cases, and technological and security constraints.

There are two primary approaches to realize the DLH vision. The first is improving the usability and flexibility of existing Big Data technologies: e.g., one could start by adding automated cluster configurations to Apache Spark. Although everyone will stand behind easier development in Spark, this approach falls short of delivering a developer experience truly aligned with the vision of the DLH, as we will discuss further below.
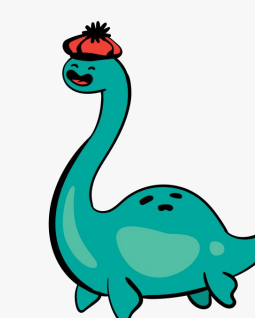
A different approach would consist in building a system from scratch based on foundational principles, while maintaining storage as a separate component; e.g., one could imagine dispensing with the Java Virtual Machine (JVM) altogether, under the assumption that the advan-
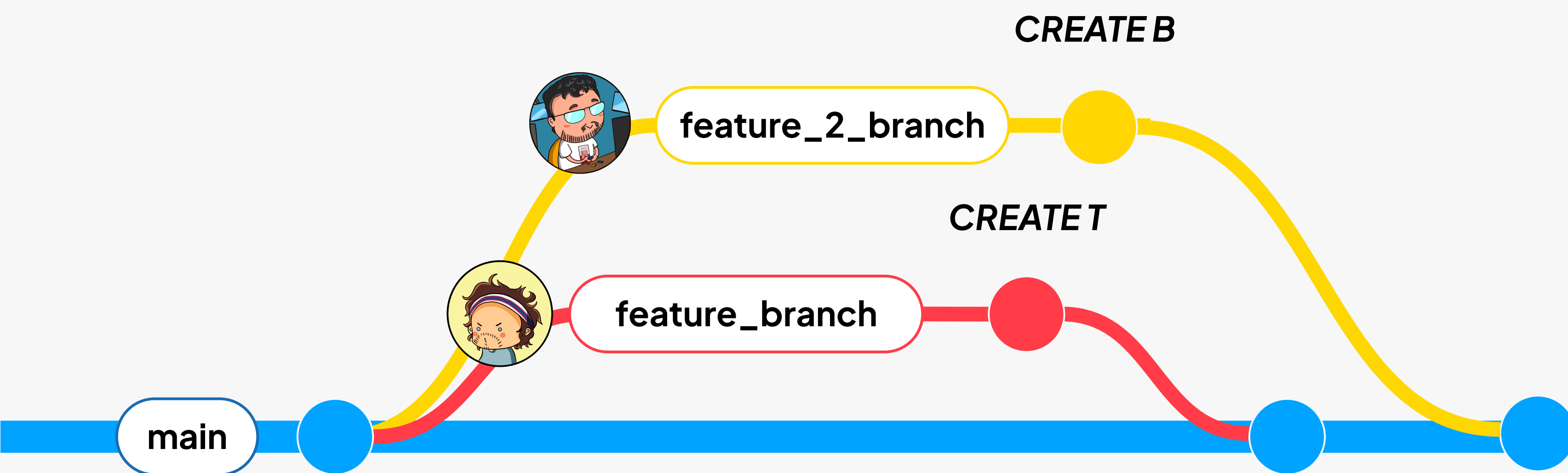
# The *composable* lakehouse

We help develop libraries we use, improve
them and often contribute back:

| Pyiceberg

| Duckdb

| Nessie

| Datafusion

| Kuzu

| Iceberg-rust
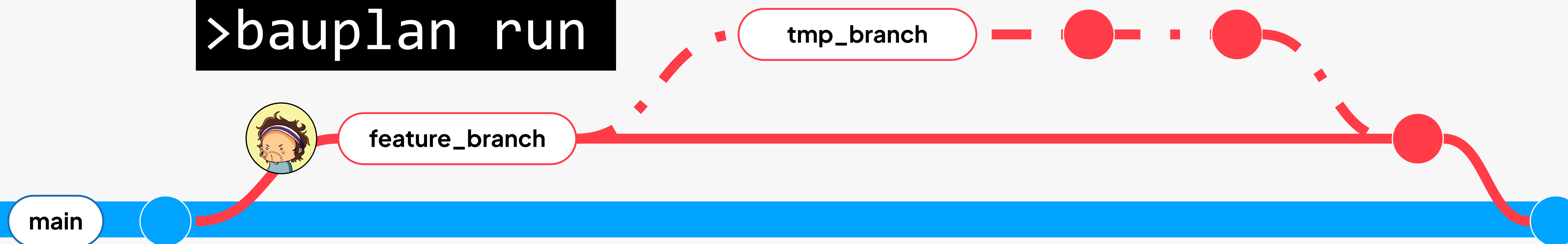
*25x
faster!*

# How much "Git" is in Git-for-data?

## Git for Data: Formal Semantics of Branching, Merging, and Rollbacks (Part 1)

How formal methods help ensure safe, reproducible workflows in data lakehouses
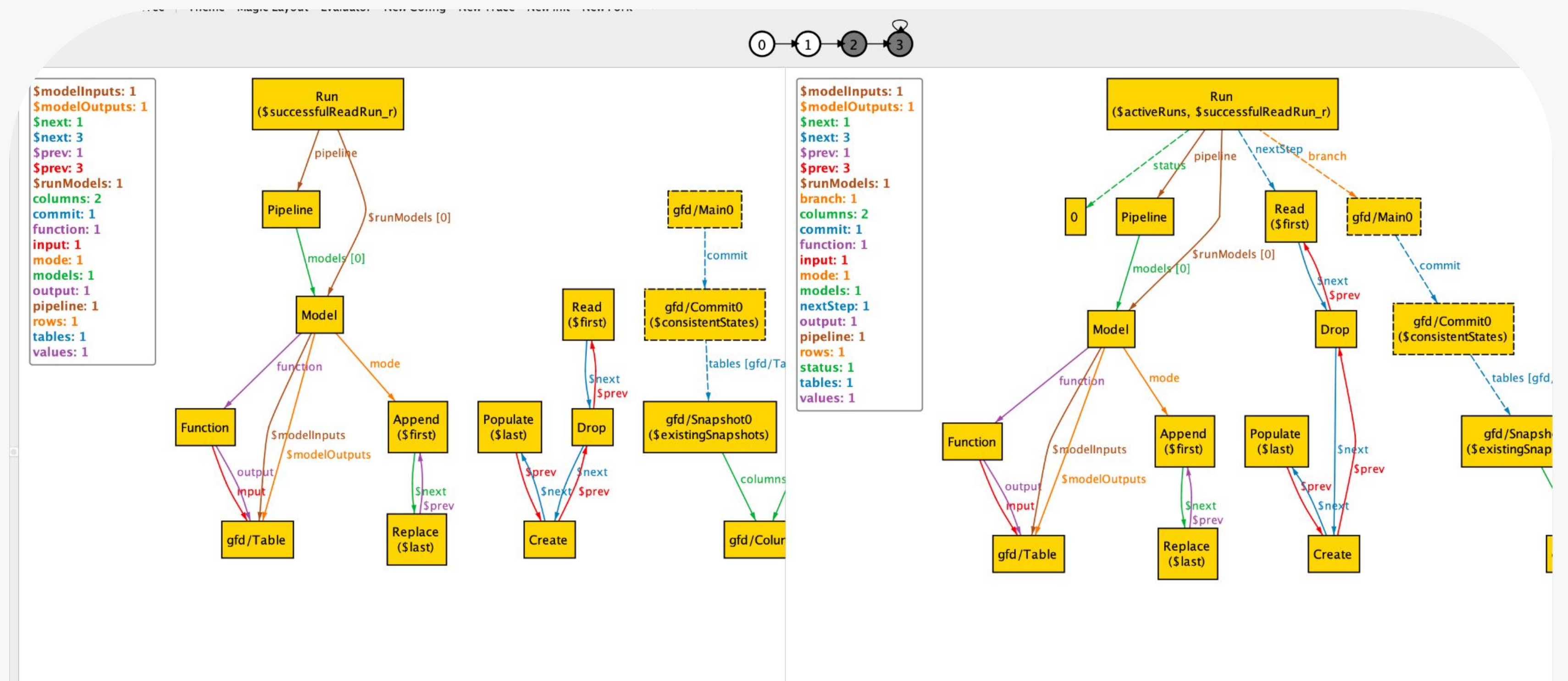
CREATE B

feature_2_branch

CREATE T

feature_branch

main

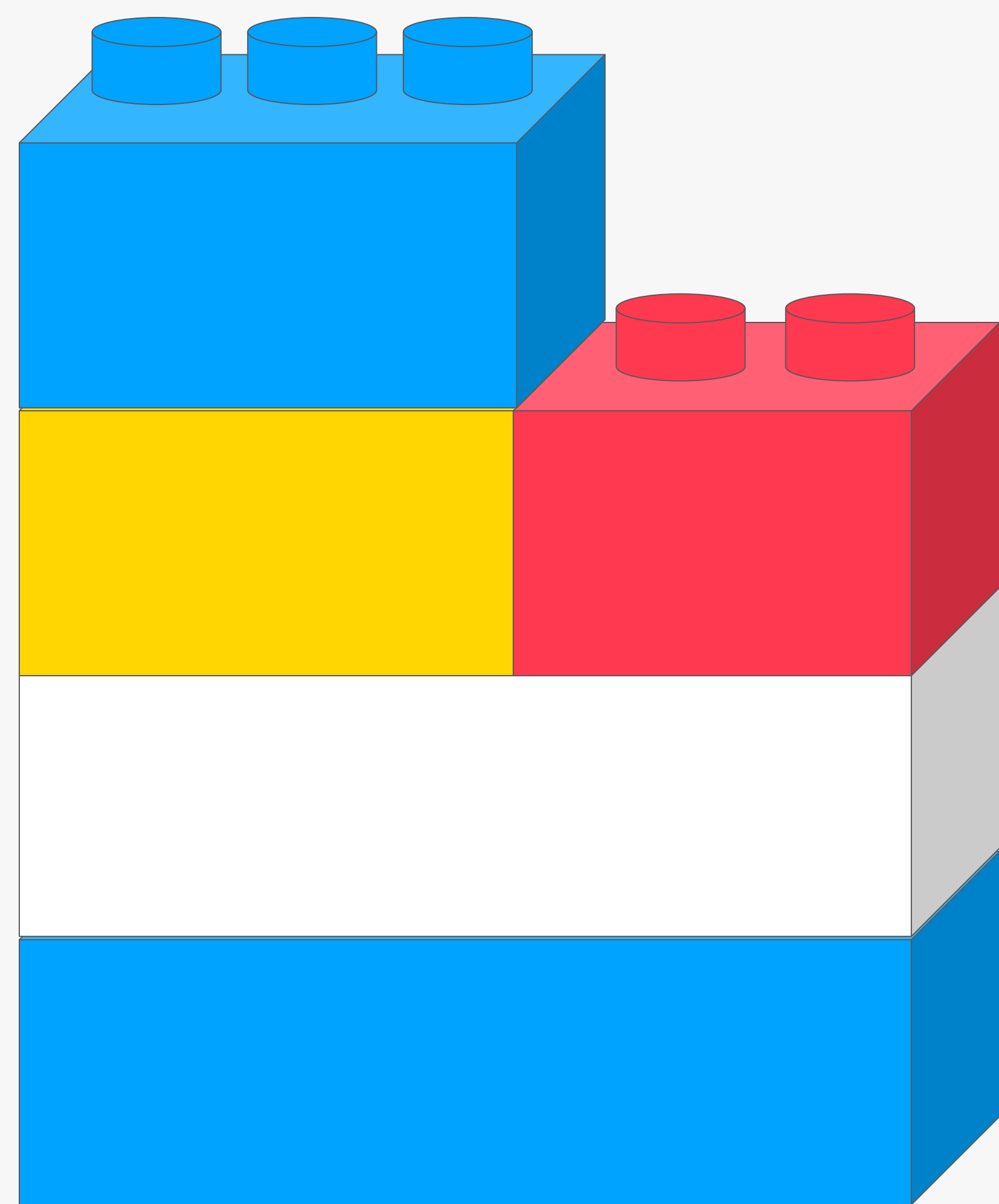# How much "database" is in Git-for-data?

```
>bauplan run
```
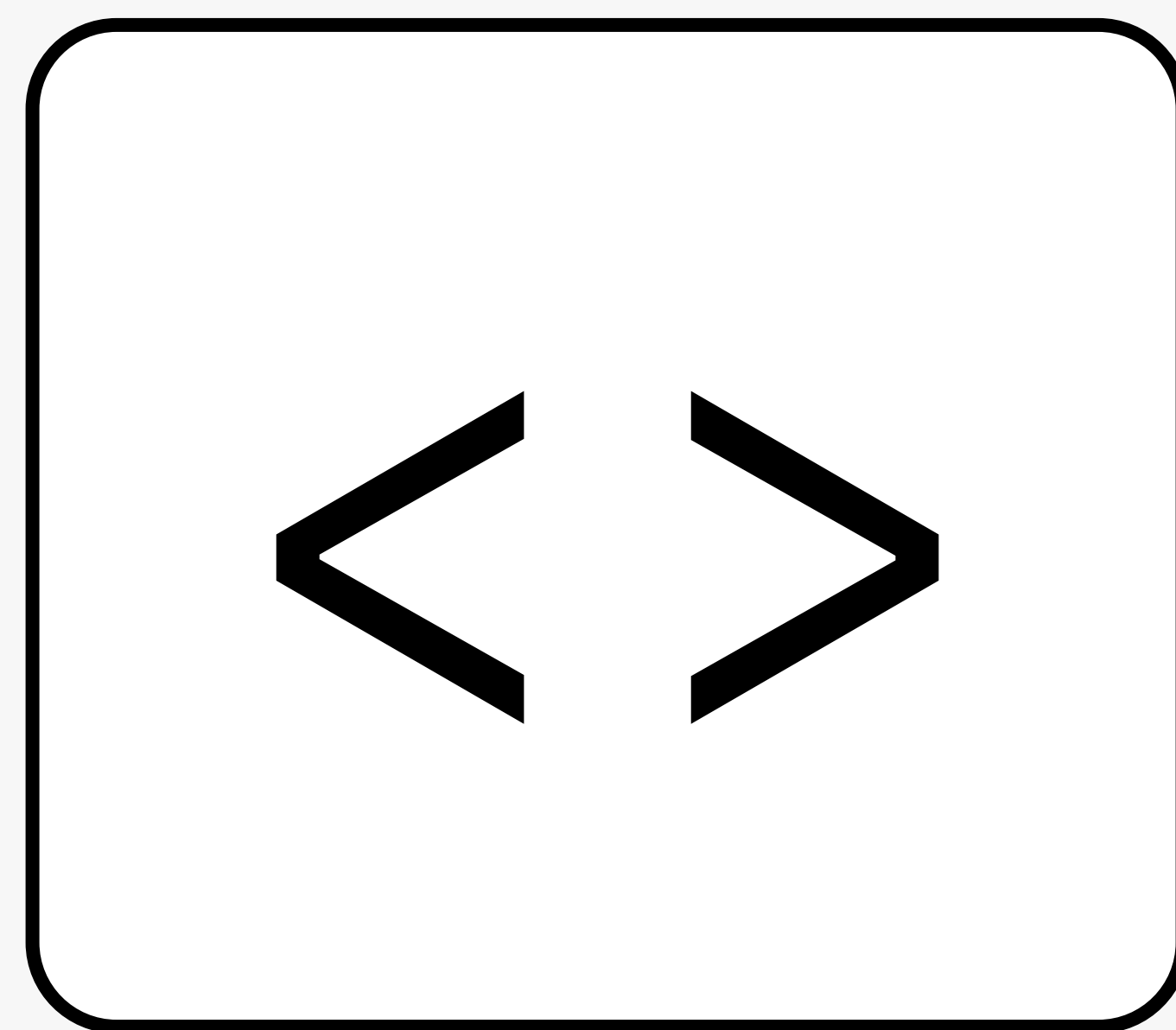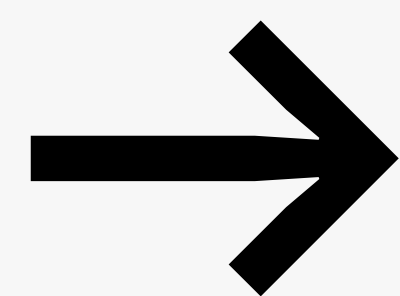
tmp_branch

feature_branch

main

*"We have discovered a truly marvelous proof of this, which this slide is too narrow to contain"*

Lightweight formal models to verify data consistency in the face of failure, and run automated checks as we add new primitives!
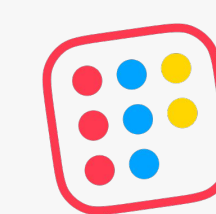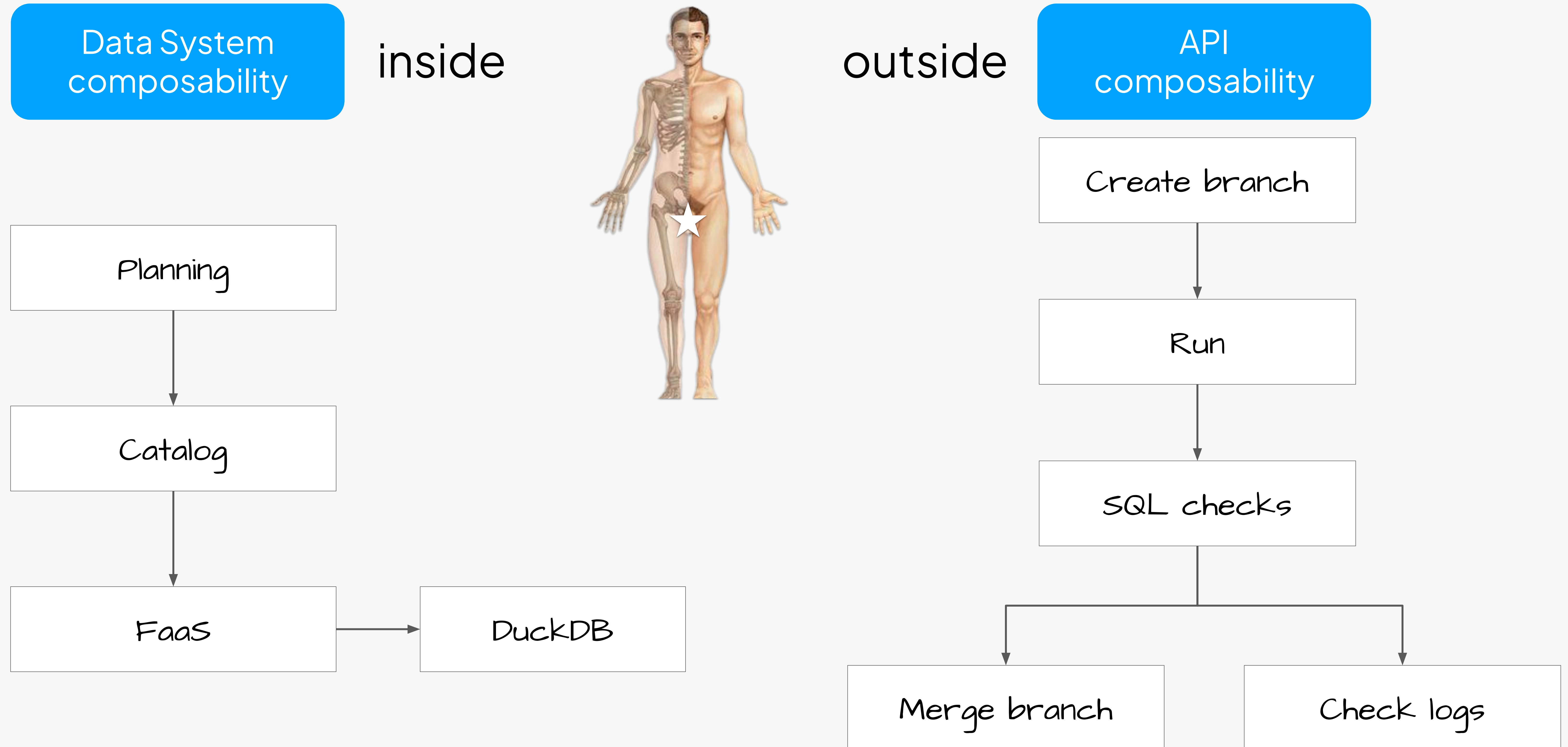
composable → programmable

# The *programmable* lakehouse

```python
1   import bauplan
2
3   @bauplan.model()
4   @bauplan.python(pip={'polars': '0.8.0'})
5   def augmented_dataset(
6       input_table='nyc_taxi',
7       columns=[col1,'col2']
8       filter="datetime ='2022-12-15",
9       model='?model'
10  ):
11      if model = 'chatgpt':
12        # init the client here
13      elif model = 'claude':
14        # init the client here
15      …
16      return predictions
17
18
19
20
```

```python
1   import bauplan
2
3   client = bauplan.Client()
4   # run models on branches
5   for i, model in enumerate(models):
6     model_branch = client.create_branch(
7         branch=f"{i}_model",
8         from_ref='main',
9         params= { "model": model },
10    )
11    run_state = client.run(
12        dir=my_pipeline,
13        branch=agent_branch
14    )
15
16  # merge the best version
17  client.merge_branch(
18        source_ref=my_best_branch,
19        into_branch='main'
        )
```
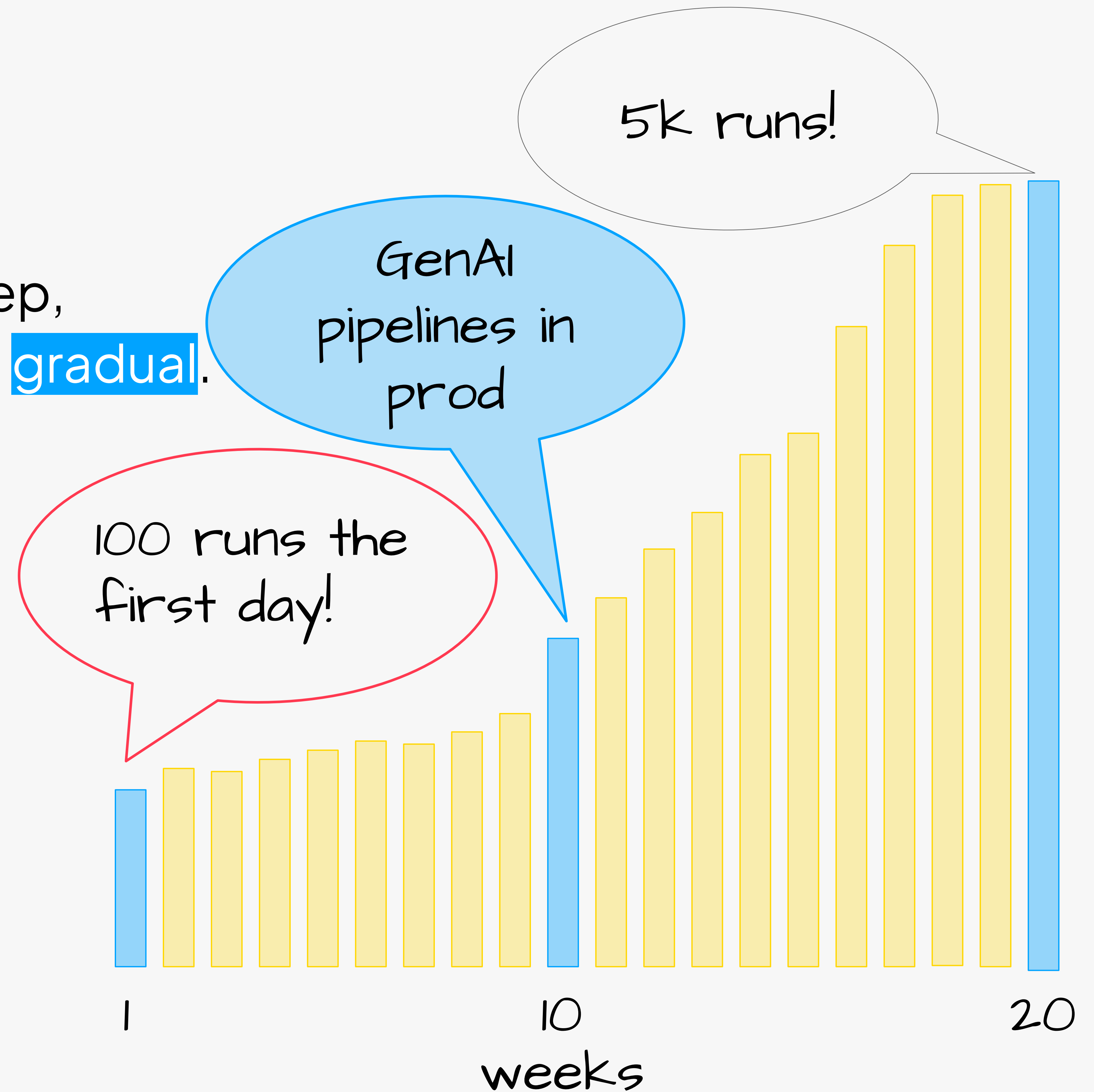
# Inside composability != outside composability

**Data System composability**    inside    outside    **API composability**

Planning

↓

Catalog

↓

FaaS → DuckDB

Create branch

↓

Run

↓

SQL checks

↓

Merge branch    Check logs

# The "API Ladder" philosophy

In order to get started, beginners need an API to be convenient.

100 runs the first day!

GenAI pipelines in prod

5k runs!

1    10    20

weeks

composable → programmable → agentic

# The *agentic* lakehouse

"**Make something idiot-proof, and someone will come up with a better idiot**"

Agents need easy-to-reason about APIs (check!), declarative infrastructure (check!) and the possibility of making mistakes without destroying downstream systems (check!).
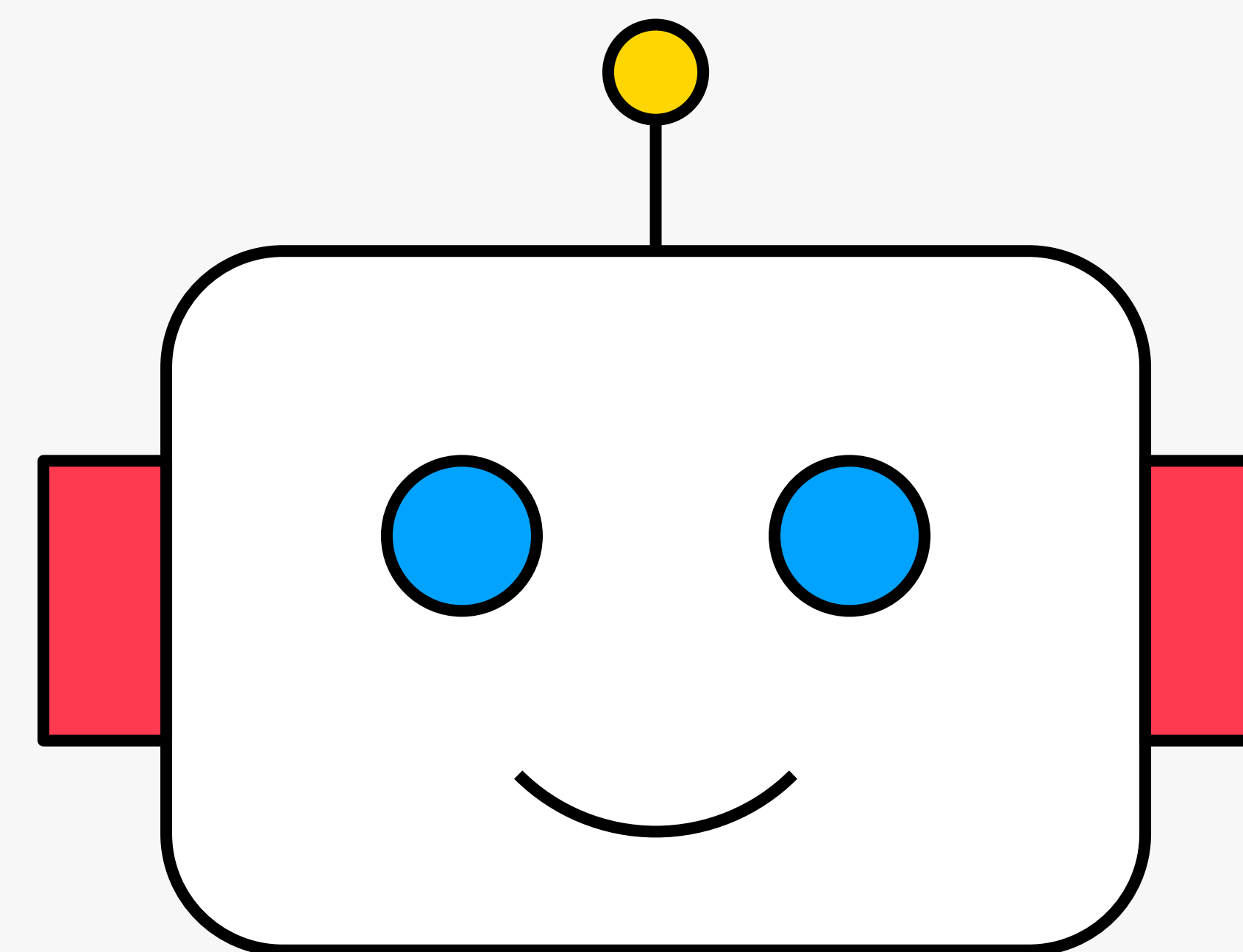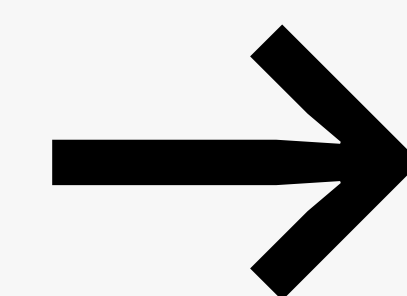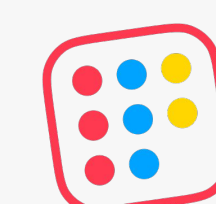
Bauplan APIs *are* the lakehouse: any model can run the full data life-cycle just with prompting!



### Safe, Untrusted, "Proof-Carrying" AI Agents: towards the agentic lakehouse

Jacopo Tagliabue
*Bauplan Labs*
NYC, USA
jacopo.tagliabue@bauplanlabs.com

Federico Bianchi
*TogetherAI*
San Francisco, USA
federico@together.ai

Ciro Greco
*Bauplan Labs*
NYC, USA
ciro.greco@bauplanlabs.com

*Abstract*—Data lakehouses manage sensitive workloads where AI automation raises risks for trust, correctness, and governance. We argue that API-first, programmable lakehouses provide the right abstractions for *safe-by-design* agentic workflows. Using Bauplan as a case study, we show how data branching and declarative environments naturally extend to agents, enabling reproducibility and observability while reducing the attack surface. We present a proof-of-concept for repairing broken data pipelines, combining Bauplan, TogetherAI, and agentic loops with correctness checks inspired by proof-carrying code. Preliminary results demonstrate both the feasibility and challenges of untrusted AI agents operating safely on production data, outlining a path towards the agentic lakehouse.

*Index Terms*—AI agents, lakehouse, data pipelines, versioning

pipelines is a canary test for agent penetration in high-stake non-trivial scenarios, which are often hard for expert humans [10], [11]. We summarize our contributions as follows:

1) we model the data pipeline life-cycle in a next-gen programmable lakehouse, Bauplan [12]: our key perspective is that traditional lakehouses resist automation because APIs are an afterthought, with no attempt to serve heterogeneous use cases with a unified interface;

2) we review common objections to automation of high-stake workloads, in the light of the proposed abstractions for repairing data pipelines: in particular, we argue that our model promotes trustworthiness and correctness both in data and code artifacts;

3) we release working code[1], showing a proof of concept for self-repairing pipelines using Bauplan, TogetherAI and an agentic loop. We share tentative results from the prototype, provide preliminary analyses

## I. INTRODUCTION

The data lakehouse is the *de facto* cloud architecture for analytics and Artificial Intelligence (AI) workloads [2], [3], thanks to storage-compute decoupling, multi-language support

# We barely scratched the surface!

# Want to know more?

**2023**

- [CDMS@VLDB 2023](#)

**2024**

- [SIGMOD 2024](#)
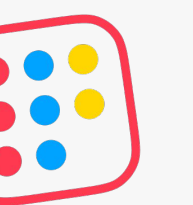- [MIDDLEWARE 2024](#) (with UMadison-Wisconsin)
- [BIG DATA 2024](#)

**2025**

- [UNDER REVIEW 2025](#) (with UMadison-Wisconsin)
- [CDMS@VLDB 2025](#) (with UChicago)

| **(Most) lakehouse** use cases can be served by a FaaS model

| **Composability** allows us to explore the design space quickly and cheaply

"It is not worth an intelligent man's time to be in the majority. By definition, there are already enough people to do that."

*G.H. Hardy*

| jacopo.tagliabue@bauplanlabs.com

| We are hiring!

bauplan